

A descriptive marker gene approach to single-cell pseudotime inference

Campbell, Kieran; Yau, Christopher

DOI:

[10.1093/bioinformatics/bty498](https://doi.org/10.1093/bioinformatics/bty498)

License:

None: All rights reserved

Document Version

Peer reviewed version

Citation for published version (Harvard):

Campbell, K & Yau, C 2018, 'A descriptive marker gene approach to single-cell pseudotime inference', *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/bty498>

[Link to publication on Research at Birmingham portal](#)

Publisher Rights Statement:

This is a pre-copyedited, author-produced version of an article accepted for publication in *Bioinformatics* following peer review. The version of record Kieran R. Campbell and Christopher Yau, 'A descriptive marker gene approach to single-cell pseudotime inference', *Bioinformatics*, bty498 is available online at: <https://doi.org/10.1093/bioinformatics/bty498>

Checked 20/06/2018

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

Subject Section

A descriptive marker gene approach to single-cell pseudotime inference

Kieran R Campbell^{1,2,5} and Christopher Yau^{2,3,4*}

¹ Department of Physiology, Anatomy and Genetics, University of Oxford, UK

² Wellcome Trust Centre for Human Genetics, University of Oxford, UK

³ Department of Statistics, University of Oxford, UK

⁴ Centre for Computational Biology, University of Birmingham, UK

⁵ Current address: Department of Statistics, University of British Columbia, Vancouver BC, Canada

*To whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Pseudotime estimation from single-cell gene expression data allows the recovery of temporal information from otherwise static profiles of individual cells. Conventional pseudotime inference methods emphasise an unsupervised transcriptome-wide approach and use retrospective analysis to evaluate the behaviour of individual genes. However, the resulting trajectories can only be understood in terms of abstract geometric structures and not in terms of interpretable models of gene behaviour.

Results: Here we introduce an orthogonal Bayesian approach termed “Ouija” that learns pseudotimes from a small set of marker genes that might ordinarily be used to retrospectively confirm the accuracy of unsupervised pseudotime algorithms. Crucially, we model these genes in terms of switch-like or transient behaviour along the trajectory, allowing us to understand why the pseudotimes have been inferred and learn informative parameters about the behaviour of each gene. Since each gene is associated with a switch or peak time the genes are effectively ordered along with the cells, allowing each part of the trajectory to be understood in terms of the behaviour of certain genes. We demonstrate that this small panel of marker genes can recover pseudotimes that are consistent with those obtained using the entire transcriptome. Furthermore, we show that our method can detect differences in the regulation timings between two genes and identify “metastable” states - discrete cell types along the continuous trajectories - that recapitulate known cell types.

Availability: An open source implementation is available as an R package at <http://www.github.com/kieranrcampbell/ouija> and as a Python/TensorFlow package at <http://www.github.com/kieranrcampbell/ouijaflow>.

Contact: kieran.campbell@stat.ubc.ca, c.yau@bham.ac.uk

Supplementary information: Supplementary text, figures, and tables are available at *Bioinformatics* online.

1 Introduction

The advent of high-throughput single-cell technologies has revolutionised single-cell biology by allowing dense molecular profiling for studies involving 100-10,000s of cells (Kalisky and Quake, 2011; Shapiro *et al.*, 2013; Macaulay and Voet, 2014; Wills and Mead, 2015). The increased

availability of single-cell data has driven the development of novel analytical methods specifically tailored to single cell properties (Stegle *et al.*, 2015; Trapnell, 2015). The difficulties in conducting genuine time-series experiments at the single-cell level has led to the development of computational techniques known as *pseudotime ordering* algorithms that extract temporal information from snapshot molecular profiles of individual cells. These algorithms exploit studies in which the captured cells behave asynchronously and therefore each is at a different stage of

some underlying temporal biological process such as cell differentiation. In sufficient numbers, it is possible to infer an ordering of the cellular profiles that correlates with actual temporal dynamics and these approaches have promoted insights into a number of time-evolving biological systems (Qiu *et al.*, 2011; Bendall *et al.*, 2014; Trapnell *et al.*, 2014; Reid and Wernisch, 2015; Hanchate *et al.*, 2015; Shin *et al.*, 2015; Haghverdi *et al.*, 2016; Setty *et al.*, 2016).

A predominant feature of current pseudotime algorithms is that they emphasise an “unsupervised” approach. The high-dimensional molecular profiles for each cell are projected on to a reduced dimensional space by using a (non)linear transformation of the molecular features. In this reduced dimensional space, it is hoped that any temporal variation is sufficiently strong to cause the cells to align against a trajectory along which pseudotime can be measured. This approach is therefore subject to a number of analysis choices including gene selection, dimensionality reduction technique, and cell ordering algorithm, all of which could lead to considerable variation in the pseudotime estimates obtained. In order to verify that any specific set of pseudotime estimates are biologically plausible, it is typical for investigators to retrospectively examine specific marker genes or proteins to confirm that the predicted (pseudo)temporal behaviour matches *a priori* beliefs. An iterative “semi-supervised” process maybe therefore be required to concentrate pseudotime algorithms on behaviours that are both consistent with the measured data and compliant with a limited amount of known gene behaviour.

2 Approach

In this paper we present an orthogonal approach implemented within a Bayesian latent variable statistical framework called ‘Ouija’ that learns pseudotimes from small panels of putative or known marker genes (Figure 1A). Our model focuses on switch-like and transient expression behaviour along pseudotime trajectories, explicitly modelling when a gene turns on or off along a trajectory or at which point its expression peaks. Crucially, this allows the pseudotime inference procedure to be understood in terms of descriptive gene regulation events along the trajectory (Figure 1B). As each gene is associated with a particular switch or peak time, it allows us to order the genes along the trajectory as well as the cells and discover which parts of the trajectory are governed by the behaviour of which genes. For example, if the pseudotimes for a set of differentiating cells run from 0 (stem cell like) to 1 (differentiated) and only two genes have switch times less than 0.25 then a researcher would conclude that the beginning of differentiation is regulated by those two genes. We further formulate a Bayesian hypothesis test as to whether a given gene is regulated before another along the pseudotemporal trajectory (Figure 1C) for all pairwise combinations of genes. Furthermore, by using such a probabilistic model we can identify discrete cell types or “metastable states” along continuous developmental trajectories (Figure 1D) that correspond to known cell types.

3 Methods

3.1 Overview

The aim of pseudotime ordering is to associate a G -dimensional expression measurement to a latent unobserved pseudotime. Mathematically we can express this as:

$$\underbrace{\mathbf{y}_n}_{\text{Expression}} = \underbrace{f}_{\text{Mapping}} \left(\underbrace{t_n}_{\text{Pseudotime}} \right) + \underbrace{\epsilon_n}_{\text{Noise}} \quad (1)$$

where the function f maps the one-dimensional pseudotime t_n for cell n to the G -dimensional observation space. The challenge lies in the

fact that both the mapping function f and the pseudotimes are unknown. Our objective here is to use parametric forms for the mapping function f that will enable relatively fast computations whilst characterising certain gene expression temporal behaviours. The specification of a statistical pseudotime algorithm therefore comes down to the choice of the mean function f and the noise model on ϵ (see Supplementary Text Section 5 for an in-depth discussion).

3.2 Input data normalisation

We index N cells by $n \in 1, \dots, N$ and G genes by $g \in 1, \dots, G$. Let $y_{ng} = [\mathbf{Y}]_{ng}$ denote the log-transformed non-negative observed cell-by-gene expression matrix. In order to make the strength parameters comparable between genes we normalise the gene expression so the approximate half-peak expression is 1 through the transformation $y_{ng} \rightarrow y'_{ng} = y_{ng}/s_g$ where s_g is a gene-specific size factor defined by

$$s_g = \frac{1}{|\mathcal{Y}_g^*|} \sum_{y_{cg}^* \in \mathcal{Y}_g^*} y_{cg}^* \quad (2)$$

$$\text{and } \mathcal{Y}_g^* = \{y_{cg} : y_{cg} > 0\}.$$

3.3 Noise model

Our statistical model can be specified as a Bayesian hierarchical model where the likelihood is given by a bimodal distribution formed from a mixture of zero-component (dropout) and an non-zero expressing cell population. If $\mu(t_n, \Theta_g)$ is the mean for cell n and gene g (evaluated at pseudotime t_n with gene-specific parameters Θ_g) then

$$\begin{aligned} \beta_0, \beta_1 &\sim \text{Normal}(0, 0.1) \\ \pi_{ng} &\sim \text{Bernoulli}(\text{logit}^{-1}(\beta_0 + \beta_1 \mu(t_n, \Theta_g))), \\ p(y_{ng} | \pi_{ng}, \mu_{ng}, \sigma_{ng}) &= \pi_{ng} \delta(y_{ng}) \\ &\quad + (1 - \pi_{ng}) T_\nu(y_{ng} | \mu(t_n, \Theta_g), \sigma_{ng}^2), \end{aligned} \quad (3)$$

where π_{ng} is the probability of observing a dropout (zero-count) in cell n gene g and T is the density function of the Student-t distribution with ν degrees of freedom.

The relationship between dropout rate and expression level is expressed as a logistic regression model (Kharchenko *et al.*, 2014). Furthermore, we impose a mean-variance relationship of the form $\sigma_{ng}^2 = (1 + \phi) \mu(t_n, \Theta_g) + \epsilon$ where ϕ is the dispersion parameter with prior $\phi \sim \text{Gamma}(\alpha_\phi, \beta_\phi)$, which is motivated by empirical observations of marker gene behaviour (Supplementary Text 4.1).

3.4 Mean functions

We then need to specify the form of the mean functions $\mu(t_n, \Theta_g)$, for which we consider both sigmoidal and transient gene behaviour. For genes we expect to be *a priori* switch-like we model

$$\mu(t_n, \Theta_g) = \frac{2\eta_g}{1 + \exp(-k_g(t_c - t_g^{(0)}))}, \quad (4)$$

where k_g and $t_g^{(0)}$ denote the activation strength and activation time parameters for each gene and η_g the average peak expression with priors $\eta_g \sim \text{Gamma}(\delta/2, 1/2)$, $k_g \sim \text{Normal}(\mu_g^{(k)}, 1/\tau_g^{(k)})$, $t_g^{(0)} \sim \text{TruncNorm}_{[0,1]}(\mu_g^{(t)}, 1/\tau_g^{(t)})$. If available, user-supplied prior information can be encoded by specifying priors on the parameters $\mu_g^{(k)}, \tau_g^{(k)}, \mu_g^{(t)}, \tau_g^{(t)}$. Otherwise, inference can be performed using uninformative hyperpriors on these parameters. Specifying $\mu_g^{(k)}$ encodes a prior belief in the strength and direction of the activation of gene g along

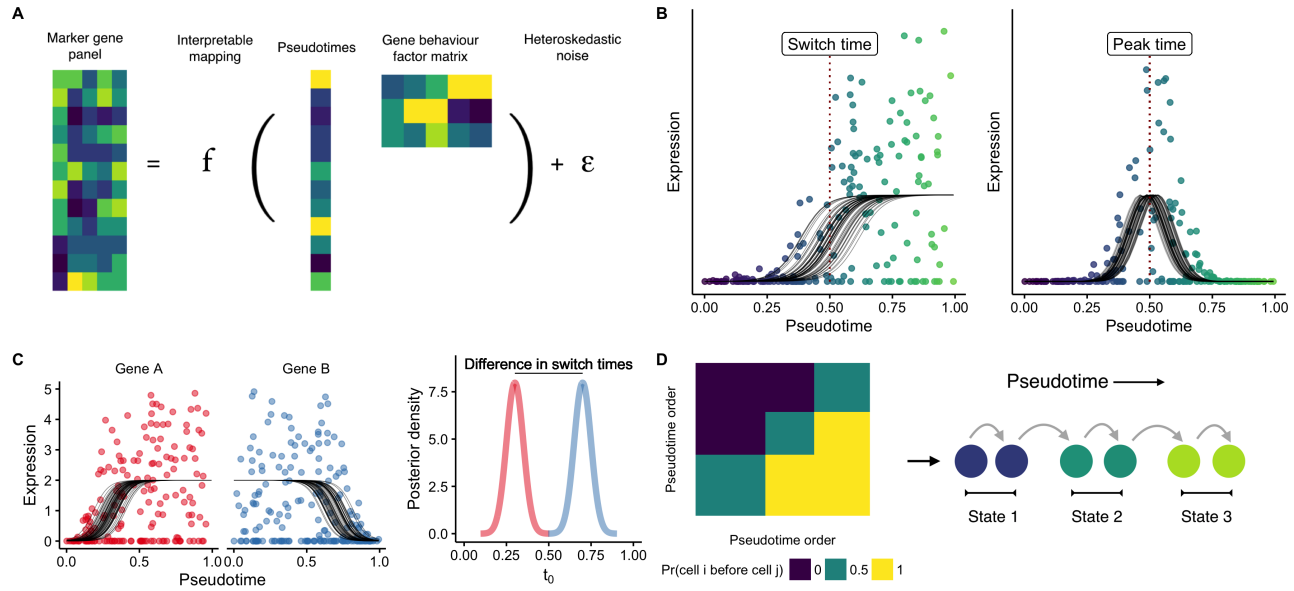


Fig. 1. Learning single-cell pseudotimes with parametric models. A Ouja infers pseudotimes using Bayesian nonlinear factor analysis by decomposing the input gene expression matrix through a parametric mapping function (sigmoidal or transient). The latent variables become the pseudotimes of the cells while the factor loading matrix is informative of different types of gene behaviour. A heteroskedastic dispersed noise model with dropout is used to accurately model scRNA-seq data. B Each gene's expression over pseudotime is modelled either as a sigmoidal shape (capturing both linear and switch-like behaviour) or through a Gaussian shape (capturing transient expression patterns). These models include several interpretable parameters including the pseudotime at which the gene is switched on and the pseudotime at which a gene peaks. C The posterior distributions over the switch and peak times can be inferred leading to a Bayesian statistical test of whether the regulation of a given gene occurs before another in the pseudotemporal trajectory. D Ouja can identify discrete cell types that exist along continuous trajectories by clustering the matrix formed by considering the empirical probability one cell is before another in pseudotime.

the trajectory with $\tau_g^{(k)}$ (inversely-) representing the strength of this belief. Similarly, specifying $\mu_g^{(t)}$ encodes a prior belief of where in the trajectory gene g exhibits behaviour (either turning on or off) with $\tau_g^{(t)}$ encoding the strength of this belief.

For the transient case we have

$$\mu(t_n, \Theta_g) = 2\eta_g \exp(-\lambda b_g(t_n - p_g)^2), \quad (5)$$

where we take $\lambda = 10$ to be a constant and with prior structure $\eta_g \sim \text{Gamma}(\delta/2, 1/2)$, $p_g \sim \text{TruncNorm}_{[0,1]}(\mu_g^{(p)}, 1/\tau_g^{(p)})$, $b_g \sim \text{TruncNorm}_{[0,\infty)}(\mu_g^{(p)}, 1/\tau_g^{(p)})$, where informative priors may be placed on p and b as before.

3.5 Inference

Under this framework learning the single-cell trajectory becomes Bayesian inference of $p(t, \Theta | \mathbf{Y})$ - the joint posterior distribution of the pseudotimes and gene behaviour parameters given the expression data. We performed posterior inference using Markov Chain Monte Carlo (MCMC) stochastic simulation algorithms, specifically the No U-Turn Hamiltonian Monte Carlo approach (Homan and Gelman, 2014) implemented in the STAN probabilistic programming language (Carpenter *et al.*, 2015). The parameter $\epsilon = 0.01$ is used to avoid numerical issues in MCMC computation. For larger marker gene panels, such as in the cell cycle analysis section, we used stochastic gradient variational Bayes implemented in Edward (Tran *et al.*, 2016) to perform approximate Bayesian inference.

4 Results & Discussion

4.1 Pseudotime inference from small marker gene panels

The transcriptomes of both single cells and bulk samples exhibit remarkable correlations across genes and transcripts. Such concerted regulation of expression is thought to be due to pathway-dependent transcription (Tegge *et al.*, 2012; Braun *et al.*, 2008) and is necessary for the field of network inference from gene expression data (Langfelder and Horvath, 2008). An example of such transcriptome wide correlations can be seen in Figure 2A for the Trapnell *et al.* (2014) dataset, where hierarchical clustering reveals a block-diagonal structure, implying an intrinsic low-dimensionality of the data that can be efficiently compressed using techniques such as principal components analysis (Supplementary Figure 1).

This redundancy of expression is often exploited by statistical models of single-cell RNA-seq data. Examples include Heimberg *et al.* (2016) where the intrinsic low-dimensionality is used to reconstruct transcriptome-wide gene expression from ultra-shallow read depths; Cleary *et al.* (2017) apply compressed sensing techniques to reconstruct high-dimensional gene expression profiles from low-dimensional random projection; and McCurdy *et al.* (2017) who propose a column subset selection procedure where a small number of genes are chosen to represent the full transcriptome. The compressibility of transcriptome data is likewise exploited by many single-cell pseudotime inference algorithms via initial dimensionality reduction steps. For example, Monocle (Trapnell *et al.*, 2014) reduces the expression data down to 2 dimensions using independent component analysis, while both TSCAN (Ji and Ji, 2016) and Waterfall (Shin *et al.*, 2015) apply PCA to reduce the data down to 2 dimensions. The implication behind such approaches that there is sufficient information in just two dimensions of the data via a linear projection to learn “transcriptome-wide” pseudotime and that the majority of expression is redundant given the low-dimensional projection.

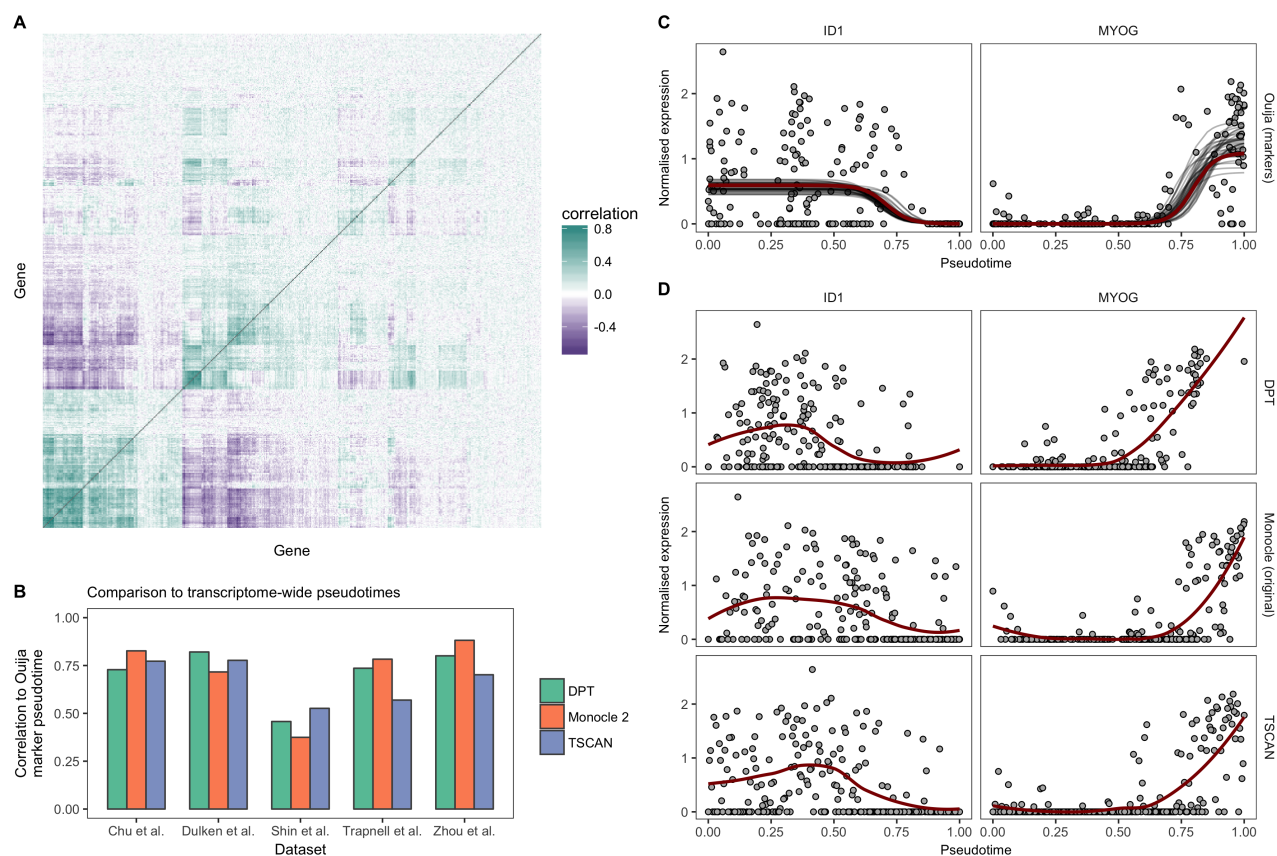


Fig. 2. Transcriptome-wide pseudotimes can be inferred from small marker gene panels. A A gene-by-gene correlation matrix for the Trapnell et al. (Trapnell et al., 2014) dataset reveals similarities in the transcriptional response of hundreds of genes. The redundancy of expression implies the information content of the transcriptome may be compressed through techniques such as principal components analysis (PCA) or by picking informative marker genes. B Comparison of pseudotimes fitted using Ouja on a small panel of marker genes to transcriptome-wide fits (using the 500 most variable genes) across five datasets using the algorithms Monocle 2, DPT, and TSCAN. The marker gene fits show high correlation to the transcriptome-wide fits with the exception of the Shin et al. (Shin et al., 2015) dataset. C Gene expression profiles for two marker genes ID1 and MYOG from the Trapnell et al. (Trapnell et al., 2014) dataset. The solid red line denotes the maximum a posteriori (MAP) Ouja fit while the grey lines show draws from the posterior mean function. D Gene expression profiles for the same genes for the algorithms DPT, Monocle 2, and TSCAN show similar expression fits, demonstrating equivalent pseudotemporal trajectories have been inferred. The solid red line denotes a LOESS fit.

In Ouja, we exploit the high gene-gene correlations by modelling a small number of *marker* genes that are representative of the whole transcriptome. Such an approach is advantageous as by modelling the data directly rather than a reduced-dimension representation we can understand the pseudotimes for each cell in terms of the behaviour of genes through time rather than abstract notions of manifolds embedded in high-dimensional space. This takes the form of a nonlinear factor analysis model, departing from previous models that have relied upon linear factor analysis (Pierson and Yau, 2015; Campbell and Yau, 2017) by introducing sigmoidal nonlinearities and transient expression functions, both of which have been successfully applied previously in post-processing of single-cell trajectories (Campbell and Yau, 2016a,b; Sander et al., 2017).

We then turn to the question of how to choose the small number of marker genes in order to fit the pseudotimes. In single-cell pseudotime studies, the cells under examination undergo a known biological process such as differentiation or cell cycle. Importantly, key marker genes associated with these processes are usually known *a priori* by investigators. These marker genes act as positive controls whose behaviour is used post-hoc to confirm the validity of the transcriptome-wide pseudotime fit. Examples include the markers of myoblast differentiation *MYH3*, *MEF2C*, and *MYOG* in Trapnell et al. (2014); the markers of neurogenesis *Gfap* and *Sox2* in Shin et al. (2015); and in Li et al. (2016) the authors tabulate the

marker genes they expect to be involved in the process along with their expected behaviour along the differentiation trajectory. Given both the widespread *a priori* knowledge of such markers and their requirement to validate transcriptome-wide pseudotime fits, we therefore propose to derive pseudotimes directly from such markers.

We first sought to test whether our model applied to small panels of marker genes could accurately recapitulate the transcriptome-wide pseudotimes inferred by popular pseudotime methods. We applied Monocle 2, DPT, and TSCAN to five publicly available single-cell RNA-seq datasets (Trapnell et al., 2014; Shin et al., 2015; Zhou et al., 2016; Dulken et al., 2017; Chu et al., 2016) using the 500 most variable genes as input (the default in packages such as Scater (McCarthy et al., 2017) for PCA representations). For each dataset, we then inferred pseudotimes using Ouja based only on a small number of marker genes reported in each paper (ranging from 5 to 12), and compared the Pearson correlation between the Ouja pseudotimes and the pseudotimes reported for each dataset (Figure 2B). There was good agreement between the marker-based pseudotimes inferred using Ouja and the transcriptome-wide pseudotimes inferred using existing algorithms, with the correlation exceeding 0.75 in the majority of comparisons.

Noting that the correlation will not be 1 unless the algorithms are identical, we sought to compare Ouja's correlation to transcriptome-wide

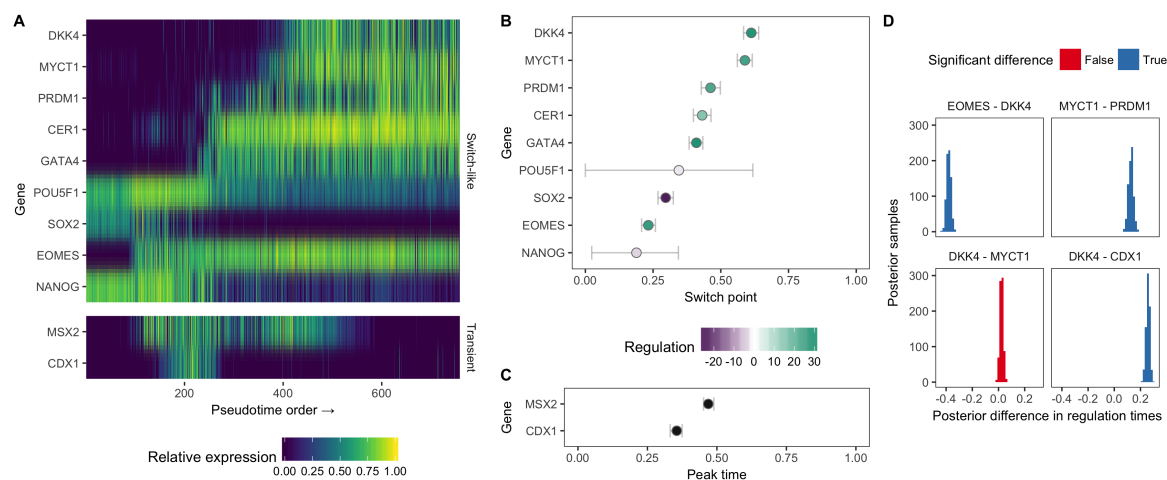


Fig. 3. Parametric models lead to pseudotimes centred around gene regulation timing. A An expression heatmap for the 9 switch-like genes and 2 transient genes in the Chu et al. dataset, with genes ordered by the posterior mean of the switch time. B-C Posterior distributions over the switch times and peak times for the 11 genes, coloured by their up or down regulation along pseudotime. The horizontal error bars show the 95% highest probability density credible intervals. D A Bayesian hypothesis test can quantify whether the posterior difference between two regulation timings (either switch or peak time) is significantly different from 0, allowing us to determine whether a given gene is regulated before or after another along pseudotime.

pseudotime to the agreement of the transcriptome-wide pseudotimes with each other. We found large variability in the agreement between existing algorithms using transcriptome-wide pseudotimes, with correlations as high as 0.93 but as low as 0.61 (Supplementary Figure 2). We found the marker-based Ouija pseudotimes have higher correlations to one of the transcriptome-wide algorithms than they have amongst each other in all but one of the datasets studied. On average, the correlation between Ouija's marker based pseudotime with the transcriptome-wide pseudotimes was around 0.1 lower than the correlation amongst the transcriptome-wide pseudotimes, though given Ouija uses around 1-2% the number of input genes we believe this is a positive result that represents transcriptome-wide pseudotimes may be inferred using interpretable, parametric models on a small number of marker genes chosen *a priori*.

This equivalence of transcriptome-wide and marker-based pseudotimes is further confirmed by examining the qualitative fit of the marker genes across the different algorithms. For example, Figure 2C shows the posterior fit of the marker-based pseudotime for two marker genes from (Trapnell et al., 2014), correctly inferring the switch-like downregulation of *ID1* and the upregulation of *MYOG*. Near identical behaviour is found using transcriptome-wide pseudotimes derived from DPT, Monocle, and TSCAN (Figure 2D). We note the low correlations of the marker-based Ouija pseudotimes with the transcriptome-wide fits for the Shin et al. dataset. Upon close inspection of the marker genes (Supplementary Figure 3) we found that the expression of four of the marker genes (*Aldoc*, *Apoe*, *Eomes*, *Sox11*) were highly correlated (the switch times are similar) whilst *Gfap* and *Stmn1* showed little variation over pseudotime. This meant that there was effectively only a single marker gene for this data set - too few for reliable marker gene-based pseudotime inference.

4.2 Gene regulation timing from marker gene-based pseudotime

Having demonstrated Ouija can accurately recapitulate transcriptome-wide pseudotimes using just small marker gene panels, we next sought to show how it allows for marker-driven interpretable inference of such trajectories. We applied Ouija to a single-cell time-series dataset of human embryonic stem cells differentiating into definitive endoderm cells (Chu et al., 2016). The authors examined the expression of key marker genes over time and found 9 to exhibit approximately switch-like behaviour

(*POU5F1*, *NANOG*, *SOX2*, *EOMES*, *CER1*, *GATA4*, *DKK4*, *MYCT1*, and *PRDM1*) with a further two exhibiting transient expression (*CDX1* and *MSX2*). We applied Ouija using noninformative priors over the behaviour parameters with no information about the capture times of the cells included.

The resulting pseudotime fit demonstrates we can understand single-cell pseudotime in terms of the behaviour of particular genes. Figure 3A shows a heatmap of the 9 switch-like genes (top) and 2 transient genes (bottom), ordered by the posterior switch time of each gene. It can be seen that the early trajectory is characterised by the expression of *NANOG*, *SOX2*, and *POU5F1*, which then leads to a cascade of switch-like activation of the remaining genes as the cells differentiate.

While transcriptome-wide pseudotime algorithms could provide similar heatmaps if the marker genes were known in advance, the key departure of Ouija is that we can quantitatively associate each gene with a region of pseudotime at which its regulation (switch time or peak time) occurs. This is illustrated in Figure 3B-C showing the posterior values for the regulation timing along with the associated uncertainty. In essence, Ouija allows us to order *genes* along trajectories as well as being able to order the cells, which provides insight into gene regulation relationships.

To approach such questions of gene regulation timings in a quantitative and rigorous manner we constructed a Bayesian hypothesis test to find out whether one gene is regulated before another given the noise in the data. If $t_{\text{Gene A}}^{(0)}$ and $t_{\text{Gene B}}^{(0)}$ are the regulation timings of genes A and B respectively, we calculate the posterior distribution $p(t_{\text{Gene A}}^{(0)} - t_{\text{Gene B}}^{(0)} | \mathbf{Y})$, and if both the lower and upper bounds of the 95% posterior credible interval fall outside 0 we say the two genes are regulated at significantly different times. We applied this to the pseudotime fit in the Chu et al. dataset, the results of which can be seen in Figure 3D for a subset of genes. The model suggests that *EOMES* is downregulated before *DKK4* and *MYCT1* is downregulated after *PRDM1*. Furthermore, it suggests the switch-like downregulation of *DKK4* occurs after the transient peak-time of *CDX1*. However, it suggests the difference in regulation timings of *DKK4* and *MYCT1* are not significantly different from zero, which could imply co-regulation.

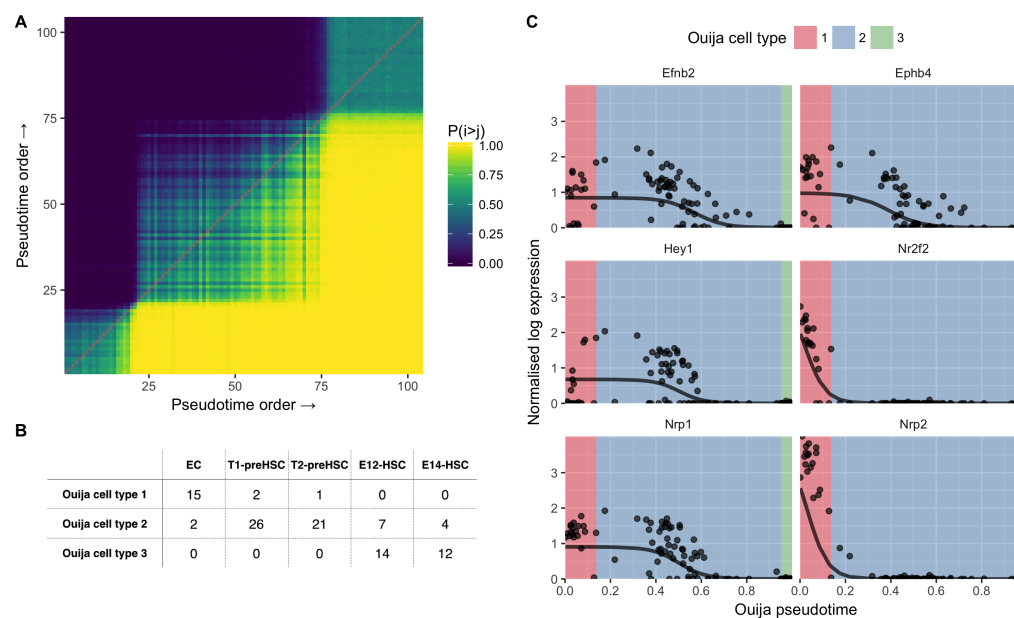


Fig. 4. Pseudotime ordering and cell type identification of haematopoietic stem cell differentiation A Consistency matrix of pseudotime ordering. Entry in the i^{th} row and j^{th} column is the proportion of times cell i was ordered before cell j in the MCMC posterior traces. Gaussian mixture modelling on the first principal component of the matrix identified three clusters that are evident in the heatmap. B Confusion matrix for cell types identified in original study (columns) and Oujia inferred (rows). Oujia inferred cluster 1 largely corresponds to EC cells, cluster 2 corresponds to pre-HSC cells while cluster 3 corresponds to HSC cells. C HSC gene expression as a function of pseudotime ordering for six marker genes. Background colour denotes the maximum likelihood estimate for the Oujia inferred cell type in that region of pseudotime.

4.3 Oujia is robust to gene behaviour misspecification

A potential disadvantage of our model is the requirement to pre-specify genes as having switch-like or transient behaviour over pseudotime, which may result in biased or erroneous pseudotimes. We noticed such an effect in the Li *et al.* (2016) dataset, where the authors pre-specified how they expected several marker genes to behave over pseudotime. Upon fitting the pseudotimes using Oujia, we noted that the genes *Mef2c* and *Pik3r2* exhibited the correct upregulation over pseudotime (Supplementary Figure 4A), but that *Scd1* that was supposed to exhibit transient, peaking expression was effectively constant along the trajectory (Supplementary Figure 4B).

We first asked whether this was a particular failing of Oujia or a result common to all pseudotime algorithms so fitted transcriptome-wide pseudotimes using TSCAN, Monocle 2 and DPT. We found remarkably low correlations between the different pseudotime algorithms (Supplementary Figure 4C), with the highest correlations reported between Oujia using markers only and Monocle 2 using the full transcriptome. Furthermore, none of the pseudotime fits displays consistent nor expected behaviour for the set of marker genes (Supplementary Figure 5).

We supplemented this with extensive simulations to discover whether Oujia is in general robust to gene behaviour misspecification. We simulated datasets where either 75% or 50% of the genes were switch-like (Supplementary Figure 4D) for 8, 12, 16 & 24 genes with 100 replications for each situation, and re-inferred the pseudotimes using Oujia assuming all genes were switch-like. The results in Supplementary Figure 4D show with 4 switch-like and 4 transient genes Oujia still achieves a median correlation greater than 0.9 with the true pseudotimes, demonstrating Oujia is highly robust to misspecification of prior knowledge of gene behaviour.

It is further possible to identify errors in the prior belief of gene behaviour without having to explicitly fit a pseudotemporal trajectory. If a dataset contains a number of switch-like and transient genes, the switch-like genes will have high absolute correlation with themselves but

low absolute correlation with the transient genes, which will in turn have high absolute correlation with themselves. This effect is exemplified in the Chu *et al.* dataset that contains 9 switch-like and 2 transient genes. A hierarchical clustering of the absolute correlations across the genes reveals the transient genes clustering separately from the switch-like genes (Supplementary Figure 6). Therefore, an investigator could corroborate their prior expectations through similar investigations.

4.4 Identifying discrete cell types along continuous developmental trajectories

We further investigated the single cell expression data from a study tracking the differentiation of embryonic precursor cells into haematopoietic stem cells (HSCs) (Zhou *et al.*, 2016). The cells begin as haemogenic endothelial cells (ECs) before successively transforming into pre-HSC and finally HSC cells. The authors identified six marker genes that would be down-regulated along the differentiation trajectory, with early down-regulation of *Nrp2* and *Nr2f2* as the cells transform from ECs into pre-HSCs, and late down-regulation of *Nrp1*, *Hey1*, *Efnb2* and *Ephb4* as the cells emerge from pre-HSCs to become HSCs. The study investigated a number of distinct cell types at different stages of differentiation: EC cells, T1 cells (*CDK45*⁻ pre-HSCs), T2 cells (*CDK45*⁺ pre-HSCs) and HSC cells at the E12 and E14 developmental stages.

We therefore sought to identify the existence of these discrete cell types along the continuous developmental trajectory. As Oujia uses a probabilistic model and inference we were able to obtain a posterior ordering “consistency” matrix (Figure 4A) where an entry in row i column j denotes the empirical probability that cell i is ordered before cell j . Performing PCA on this matrix gives a rank-one representation of cell-cell continuity, which is then clustered using a Gaussian mixture model to find discrete cell states along the continuous trajectory (where the number of states is chosen such that the Bayesian information criterion is maximised).

Applying this methodology to the Zhou *et al.* dataset uncovered three metastable groups of cells corresponding to endothelial, pre-HSCs and HSCs respectively (Figure 4B). Misclassifications within cell types (T1/T2 and E12/E14 cells) could be explained by examining a principal components analysis of the global expression profiles (Supplementary Figure 7) which suggests that these cell types are not completely distinct in terms of expression. When examining the inferred pseudotime progression of each marker gene (Figure 4C), these three metastable states corresponded to the activation of all genes at the beginning of pseudotime time, the complete inactivation of all the marker genes at the end of the pseudotime and an intervening transitory period as each marker gene turns off. Each metastable state clearly associates with a particular cell type with *Nrp2* and *Nr2f2* exhibiting early down-regulation and *Nrp1*, *Hey1*, *Efnb2* and *Ephb4* all exhibiting late down-regulation. Using this HSC formation system as a proof-of-principle it is evident that, if a small number of switch-like marker genes are known, it is possible to recover signatures of temporal progression using Ouija and that these trajectories are compatible with real biology.

To show the widespread applicability of this method we applied it to two further publically available datasets. Dulken *et al.* (Dulken *et al.*, 2017) examined the trajectory of quiescent neural stem cells (qNSCs) as they differentiate into activated neural stem cells (aNSCs) and neural progenitor cells (NPCs). Applying Ouija's clustering-along-pseudotime revealed seven distinct clusters (Supplementary Figure 8; Supplementary Table 1) with clusters 1-2 corresponding to early and late qNSCs, cluster 3 defining the qNSC to aNSC transition, clusters 4-6 corresponding to early to late aNSCs and cluster 7 defining the aNSC to NPC transition. We similarly applied this method to the Chu *et al.* dataset of time-series scRNA-seq that identified 8 distinct clusters along pseudotime (Supplementary Figure 9; Supplementary Table 2). Clusters 1-4 track the cells as the progress through the 4 stages from 0 hours to 36 hours, while clusters 5-8 track the 3 stages from 36 hours to 96h hours but with much more heterogeneity within each cluster, which is expected due to the longer time-scales considered.

4.5 Scalable pseudotime inference using TensorFlow

Finally, we wanted to consider a study composed of a large panel of putative marker genes to determine if Ouija could automatically identify genes satisfying its behavioural constraints. We identified a single-cell RNA-seq study (Kowalczyk *et al.*, 2015) that examined variation between individual hematopoietic stem and progenitor cells from two mouse strains (C57BL/6 and DBA/2) as they age. Principal component analysis for each cell type and age showed a striking association of the top principal components with cell cycle-related genes (Figure 5A), indicating that transcriptional heterogeneity was dominated by cell cycle status. They scored each cell for its likely cell cycle phase using signatures based on functional annotations (of the Gene Ontology Consortium *et al.*, 2009) and profiles from synchronized HeLa cells (Whitfield *et al.*, 2002) for the G1/S, S, G2, and G2/M phases.

We investigated if Ouija could be used to identify cell cycle phase, treating the inferential problem as a continuous pseudotime process and assuming all genes as candidate switch genes. We applied Ouija to 1,008 C57BL/6 HSCs using 374 GO cell cycle genes that satisfied gene selection criteria used in the original study. This large number of genes and cells makes inference using Hamiltonian Monte Carlo (HMC) slow so we implemented a second version of Ouija (termed *Ouijaflow*) using the probabilistic programming language Edward (Tran *et al.*, 2016) based on TensorFlow (Abadi *et al.*, 2016). This performs fast approximate Bayesian inference using stochastic gradient variational inference (Supplementary Figure 10).

The estimated pseudotime progression given by Ouija recapitulates the trajectory observed in principal component space (Figure 5A). The estimated pseudotime distribution correlates well with the cell cycle phase categorisation given in the original study (Figure 5C). Furthermore, we identified 88 genes with large activation strengths indicating strong switching-on behaviour (Figure 5D). Ordering the genes by activation time demonstrates a cascade of expression activation across these 88 genes over cell cycle progression with the quiescent (G_0) indicated by complete inactivation of all 88 genes (Figure 5E,F). The explicit parametric model assumed by Ouija makes this gene selection and ordering process simple and quantitative compared to a non-parametric approach that would require some retrospective analysis or visual inspection.

5 Conclusion

We have developed a novel approach for pseudotime estimation based on modelling switch-like and transient expression behaviour for a small panel of marker genes chosen *a priori*. Our strategy provides an orthogonal and complementary approach to unsupervised whole-transcriptome methods that do not explicitly model any gene-specific behaviours and do not readily permit the inclusion of prior knowledge.

We demonstrate that the selection of a few marker genes allows comparable pseudotime estimates to whole transcriptome methods on real single cell data sets. Furthermore, using a parametric gene behaviour model and full Bayesian inference we are able to recover posterior uncertainty information about key parameters, such as the gene activation time, allowing us to explicitly determine a potential ordering of gene (de)activation and peaking events over pseudotime. The posterior ordering uncertainty can also be used to identify homogeneous metastable phases of transcriptional activity that might correspond to transient, but discrete, cell states. In summary, Ouija provides a novel contribution to the increasing plethora of pseudotime estimation methods available for single cell gene expression data.

Funding

KRC is supported by a doctoral studentship from the UK Medical Research Council and a postdoctoral fellowship from the Canadian Statistical Sciences Institute. CY is supported by the UK Medical Research Council.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2016). TensorFlow: Large-Scale machine learning on heterogeneous distributed systems.
- Bendall, S. C., Davis, K. L., Amir, E.-A. D., Tadmor, M. D., Simonds, E. F., Chen, T. J., Shenfeld, D. K., Nolan, G. P., and Pe'er, D. (2014). Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. *Cell*, **157**(3), 714–25.
- Braun, R., Cope, L., and Parmigiani, G. (2008). Identifying differential correlation in gene/pathway combinations. *BMC bioinformatics*, **9**(1), 488.
- Campbell, K. R. and Yau, C. (2016a). Order under uncertainty: Robust differential expression analysis using probabilistic models for pseudotime inference. *PLoS Comput. Biol.*, **12**(11), e1005212.
- Campbell, K. R. and Yau, C. (2016b). switchde: inference of switch-like differential expression along single-cell trajectories. *Bioinformatics*.
- Campbell, K. R. and Yau, C. (2017). Probabilistic modeling of bifurcations in single-cell gene expression data using a bayesian mixture of factor analyzers. *Wellcome Open Res*, **2**, 19.

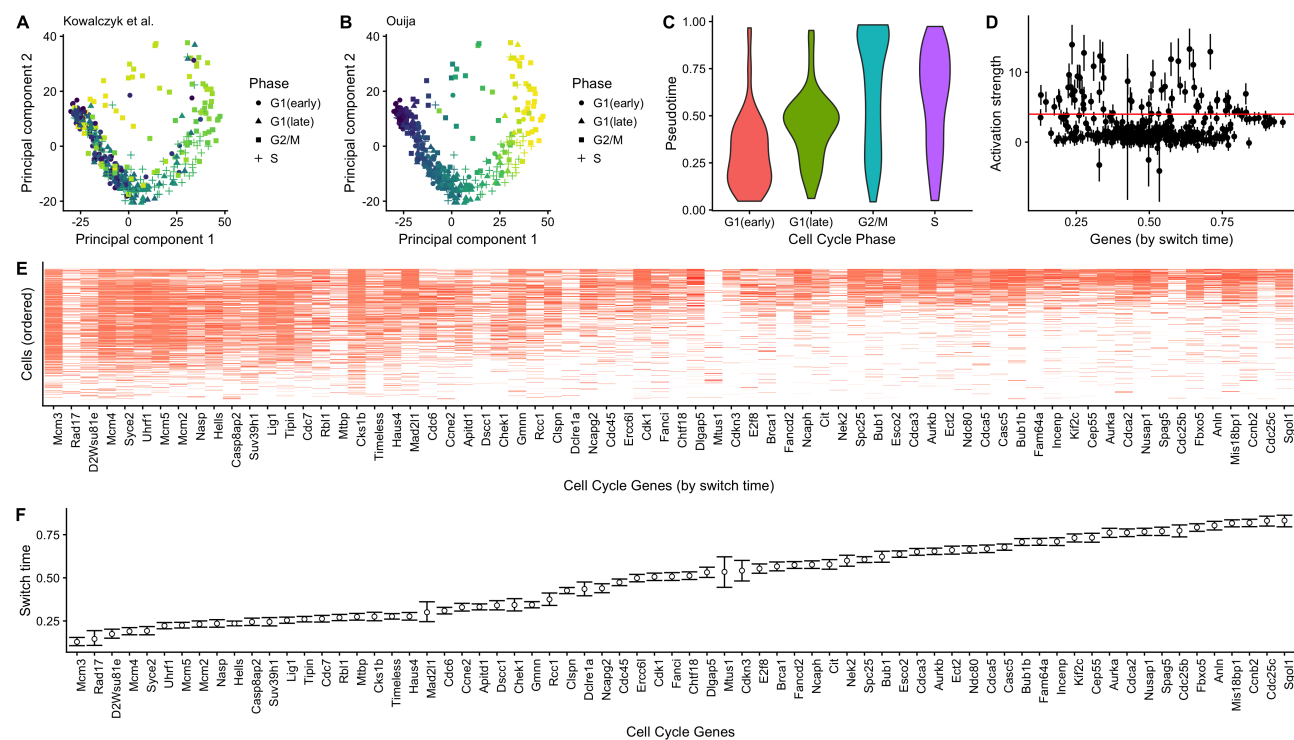


Fig. 5. Cell cycle phase prediction. Principal component representation of hematopoietic stem cells coloured according to (A) the original cell cycle progression score (Kowalczyk et al., 2015) and (B) Oujia - cell cycle classes indicated are based on original study classifications. (C) Distribution of Oujia inferred pseudotime versus the original cell cycle classifications. (D) Estimated activation strengths for the 374 cell cycle gene panels. (E) Gene expression profile for 88 switch-like genes with cells ordered by pseudotime and (F) genes ordered by activation time.

- Carpenter, B., Gelman, A., Hoffman, M., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P., and Riddell, A. (2015). Stan: a probabilistic programming language. *Journal of Statistical Software*.
- Chu, L.-F., Leng, N., Zhang, J., Hou, Z., Mamott, D., Vereide, D. T., Choi, J., Kendzior, C., Stewart, R., and Thomson, J. A. (2016). Single-cell rna-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm. *Genome biology*, **17**(1), 173.
- Cleary, B., Cong, L., Lander, E., and Regev, A. (2017). Composite measurements and molecular compressed sensing for highly efficient transcriptomics. *bioRxiv*, page 091926.
- Dulken, B. W., Leeman, D. S., Boutet, S. C., Hebestreit, K., and Brunet, A. (2017). Single-cell transcriptomic analysis defines heterogeneity and transcriptional dynamics in the adult neural stem cell lineage. *Cell reports*, **18**(3), 777–790.
- Haghverdi, L., Büttner, M., Wolf, F. A., Büttner, F., and Theis, F. J. (2016). Diffusion pseudotime robustly reconstructs lineage branching. *Nature Methods*.
- Hanchate, N. K., Kondoh, K., Lu, Z., Kuang, D., Ye, X., Qiu, X., Pachter, L., Trapnell, C., and Buck, L. B. (2015). Single-cell transcriptomics reveals receptor transformations during olfactory neurogenesis. *Science*, pages science.aad2456–.
- Heimberg, G., Bhatnagar, R., El-Samad, H., and Thomson, M. (2016). Low dimensionality in gene expression data enables the accurate extraction of transcriptional programs from shallow sequencing. *Cell systems*, **2**(4), 239–250.
- Homan, M. D. and Gelman, A. (2014). The no-u-turn sampler: Adaptively setting path lengths in hamiltonian monte carlo. *The Journal of Machine Learning Research*, **15**(1), 1593–1623.
- Ji, Z. and Ji, H. (2016). Tscan: Pseudo-time reconstruction and evaluation in single-cell rna-seq analysis. *Nucleic acids research*, page gkw430.
- Kalisky, T. and Quake, S. R. (2011). Single-cell genomics. *Nature methods*, **8**(4), 311–314.
- Kharchenko, P. V., Silberstein, L., and Scadden, D. T. (2014). Bayesian approach to single-cell differential expression analysis. *Nature methods*, **11**(7), 740–2.
- Kowalczyk, M. S., Tirosh, I., Heckl, D., Rao, T. N., Dixit, A., Haas, B. J., Schneider, R. K., Wagers, A. J., Ebert, B. L., and Regev, A. (2015). Single-cell rna-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome research*, **25**(12), 1860–1872.
- Langfelder, P. and Horvath, S. (2008). Wgcna: an r package for weighted correlation network analysis. *BMC bioinformatics*, **9**(1), 559.
- Li, J., Luo, H., Wang, R., Lang, J., Zhu, S., Zhang, Z., Fang, J., Qu, K., Lin, Y., Long, H., et al. (2016). Systematic reconstruction of molecular cascades regulating gp development using single-cell rna-seq. *Cell reports*, **15**(7), 1467–1480.
- Macaulay, I. C. and Voet, T. (2014). Single cell genomics: advances and future perspectives. *PLoS genetics*, **10**(1), e1004126.
- McCarthy, D. J., Campbell, K. R., Lun, A. T. L., and Wills, Q. F. (2017). Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics*.
- McCurdy, S., Ntranos, V., and Pachter, L. (2017). Column subset selection for single-cell rna-seq clustering. *bioRxiv*, page 159079.
- of the Gene Ontology Consortium, R. G. G. et al. (2009). The gene ontology's reference genome project: a unified framework for functional annotation across species. *PLoS Comput Biol*, **5**(7), e1000431.
- Pierson, E. and Yau, C. (2015). ZIFA: Dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome biology*, **16**(1), 241.
- Qiu, P., Simonds, E. F., Bendall, S. C., Gibbs Jr, K. D., Bruggner, R. V., Linderman, M. D., Sachs, K., Nolan, G. P., and Plevritis, S. K. (2011). Extracting a cellular hierarchy from high-dimensional cytometry data with spade. *Nature biotechnology*, **29**(10), 886–891.
- Reid, J. E. and Wernisch, L. (2015). Pseudotime estimation: deconfounding single cell time series. *bioRxiv*, page 019588.
- Sander, J., Schultze, J. L., and Yosef, N. (2017). ImpulSED: detection of differentially expressed genes in time series data using impulse models. *Bioinformatics*, **33**(5), 757–759.
- Setty, M., Tadmor, M. D., Reich-Zeliger, S., Angel, O., Salame, T. M., Kathail, P., Choi, K., Bendall, S., Friedman, N., and Pe'er, D. (2016). Wishbone identifies bifurcating developmental trajectories from single-cell data. *Nature biotechnology*, **34**(6), 637–645.
- Shapiro, E., Biezuner, T., and Linnarsson, S. (2013). Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nature reviews. Genetics*, **14**(9), 618–30.
- Shin, J., Berg, D. A., Zhu, Y., Shin, J. Y., Song, J., Bonaguidi, M. A., Enikolopov, G., Nauen, D. W., Christian, K. M., Ming, G.-I., and Song, H. (2015). Single-Cell RNA-Seq with Waterfall Reveals Molecular Cascades underlying Adult Neurogenesis. *Cell Stem Cell*, **17**(3), 360–372.

- Stegle, O., Teichmann, S. a., and Marioni, J. C. (2015). Computational and analytical challenges in single-cell transcriptomics. *Nature Reviews Genetics*, **16**(3), 133–145.
- Tegge, A. N., Caldwell, C. W., and Xu, D. (2012). Pathway correlation profile of gene-gene co-expression for identifying pathway perturbation. *PloS one*, **7**(12), e52127.
- Tran, D., Kucukelbir, A., Dieng, A. B., Rudolph, M., Liang, D., and Blei, D. M. (2016). Edward: A library for probabilistic modeling, inference, and criticism.
- Trapnell, C. (2015). Defining cell types and states with single-cell genomics. *Genome Res*, **25**(10), 1491–8.
- Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., Lennon, N. J., Livak, K. J., Mikkelsen, T. S., and Rinn, J. L. (2014). The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology*, **32**(4), 381–6.
- Whitfield, M. L., Sherlock, G., Saldanha, A. J., Murray, J. I., Ball, C. A., Alexander, K. E., Matese, J. C., Perou, C. M., Hurt, M. M., Brown, P. O., *et al.* (2002). Identification of genes periodically expressed in the human cell cycle and their expression in tumors. *Molecular biology of the cell*, **13**(6), 1977–2000.
- Wills, Q. F. and Mead, A. J. (2015). Application of single cell genomics in cancer: Promise and challenges. *Human molecular genetics*, page ddv235.
- Zhou, F., Li, X., Wang, W., Zhu, P., Zhou, J., He, W., Ding, M., Xiong, F., Zheng, X., Li, Z., *et al.* (2016). Tracing haematopoietic stem cell formation at single-cell resolution. *Nature*, **533**(7604), 487–492.